



## Mutual Aid: An Indirect Evolution Analysis

Tarik Tazdaït, Alejandro Caparros, Jean-Chrsitophe Péreau

### ► To cite this version:

Tarik Tazdaït, Alejandro Caparros, Jean-Chrsitophe Péreau. Mutual Aid: An Indirect Evolution Analysis. 2008. halshs-00275386

**HAL Id: halshs-00275386**

**<https://shs.hal.science/halshs-00275386>**

Preprint submitted on 23 Apr 2008

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## **Mutual Aid: an Indirect Evolution Analysis**

**Alejandro CAPARRÓS**

Spanish Council for Scientific Research (CSIC), Institute for Public Goods and Policies (IPP).  
Albasanz 26, 3E4. 28037 Madrid.

**Jean-Christophe PEREAU**

University of Marne-la-Vallée - O.E.P, Cité Descartes, 5 boulevard Descartes, Champs Sur  
Marne, 77454 Marne-la-Vallée cedex 2, France. E-mail: pereau@univ-mlv.fr. Tel:  
+3360957058. Fax: +33160957050. And C.N.R.S - E.H.E.S.S - CIRED

**Tarik TAZDAÏT**

C.N.R.S - E.H.E.S.S - CIRED, Jardin Tropical - 45 bis, avenue de la Belle Gabrielle, 94736  
Nogent Sur Marne cedex, France.

### **Abstract.**

This paper studies the concept of “mutual aid” developed by Kropotkin, which implies cooperation as a strategic choice. We study this concept in a Sequential Prisoners’ Dilemma in a non-cooperative framework and in an indirect evolution framework (with complete and incomplete information). We systematically compare this game with one that models Kant’s moral. In the non-cooperative framework both moral concepts imply multiple equilibria. In the indirect evolution framework with complete information Kropotkin’s moral concept leads to generalized cooperation, while Kant’s rules lead towards general defection. In the indirect evolution framework with incomplete information both moral approaches favor selfishness. However, if some agents have an imperfect detection technology cooperative behavior will not disappear in Kropotkin’s case, while it will vanish with Kant’s morality.

**Key words:** mutual aid, non-cooperative game theory, indirect evolution, Kropotkin, Kant.

## 1. Introduction

Many authors have argued that morality could solve problems for which individually rational actions lead to ineffective collective solutions, as in the prisoner's dilemma (PD). Morality is analyzed in economics mainly focusing on philosophical references, such as Kant's categorical imperative (Laffont, 1975; Etzioni, 1987; Bordinon, 1990, Bergstrom, 1995, Roemer, 1996, among others). An alternative approach is proposed by Becker (1974, 1981), who explains cooperation by supposing altruistic individuals. However, Becker also treats altruism as a moral issue, although not explicitly. In any case, in these approaches motivations and behaviors merge (Ballet, 2000) and the analysis focuses on behavior, forsaking connections between motivations and behavior.

On the contrary, the concept of cooperation developed by Kropotkin (1902) in "Mutual Aid: a Factor of Evolution" has not deserved the attention that it probably merits (Fong, Bowles and Gintis, 2006). This author studies the importance of mutual-aid in human evolution. Analyzing Siberian tribes, Polynesian islanders, medieval corporations and incipient industrial societies, Kropotkin stresses the importance of the mutual-aid principle. Since the Stone Age, mutual-aid was essential for survival and progress, leading men to live in tribes and clans. Men soon developed a great number of social institutions which determined the main features of progress: collective hunting, collective defense or collective possession of territories. As new needs came up, villages, cities and finally States appeared. For Kropotkin, the mutual aid principle, transmitted as a heritage of a very long evolution, is strong enough to survive even under the most authoritative forms of State. Although Kropotkin by no means denies the importance of the "individual ego", for him, progress results from the capacity of the mutual aid principle to counter-balance the "individual ego".

In Kropotkin's analysis cooperation is a strategic choice, and not a direct result of moral considerations (as it is with Kant's rules). The agent arbitrates between individual interest (selfishness) and cooperation (the mutual aid principle). Depending on the context, this may lead to cooperation. More precisely, in certain economic situations, individuals will widen their space of strategies, including elements able to promote cooperation, to overcome dead-ends created by individualistic behaviors.

A similar idea can be found in Margolis (1982). According to this author, individuals act as if their interest would be referred to larger entities than themselves. They appear to act, not according to their own interest, but according to the interest of the community. He adds that only under this condition society can emerge. A similar idea can be found in the "we-intentions" defined in Tuomela and Miller (1988) or in the "joint goals" in Tuomela (1990).

Thus, for some authors, agents are not only guided by selfishness. They are supposed to arbitrate between their personal interest and a more general benefit. Of course, this implies a particular form of preferences, not only restricted to individual interests, but including a collective aspect. To study this form of cooperation we analyze Kropotkin's mutual aid concept as a "feeling"<sup>1</sup> favoring cooperation in a sequential prisoner's dilemma (SPD). We study this SPD in an indirect evolution framework in which cooperative equilibria can emerge via evolutionary rather than strategic considerations. The indirect evolution approach allows endogenous derivation of preferences (Güth and Yaari, 1992; Güth and Kliemt, 1992; Köningstein and Müller, 2000; Berninghaus *et al.*, 2003; Güth and Napel, 2006), and can be seen as a generalization of rational choice models, which usually take preferences as given. In this approach, short-run decisions are taken rationally (*i.e.* maximizing a utility function that incorporates monetary and non-monetary elements) while evolutionary success is exclusively given by the monetary payoff<sup>2</sup>. Thus, preferences for monetary payoffs are stable (including the standard assumption that more is preferred to less), although moral preferences (related to the non-monetary part of the utility functions) may vary (for a general analysis of the impact of incomplete information in these type of models see Ok and Vega-Redondo (2001) and Dekel, Ely and Yilankaya (2007)). It is not easy to evaluate the real-life pertinence of this assumption, since moral preferences are difficult to observe. Nevertheless, an interesting parallel can be made with the choice of religion. Each religion imposes a given set of moral principles which may impact our utility function. However, if we observe that people of other religions systematically perform better in monetary terms (or have a better social consideration) we may change our religion. In historical terms, fluctuations in dominant religions in a given country have been constant and, not always, associated with violent imposition. Thus, the assumption that changes in religion (or in other set of moral rules) are driven by imitation (*i.e.* following an evolutionary pattern) is reasonable. Of course, the assumption that agents maximize their utility (incorporating elements imposed by their religion or moral code) in the short run is also perfectly reasonable.

---

<sup>1</sup> In "Mutual Aid, a Factor of Evolution" Kropotkin (1902) uses the term "mutual aid" in opposition to selfishness. He uses, among other expressions, the mutual aid "tendency", "feeling" or "instinct".

<sup>2</sup> It is not essential to the indirect evolution approach to impose monetary payoff as the indicator of evolutionary success, it is enough to consider that short-run strategic decisions are taken maximizing utility and that evolutionary success is given by a different measure. However, most applications have used monetary success as the measure of evolutionary success (Köningstein and Müller, 2000). Nevertheless, "social success" could also be seen as the measure of evolutionary success.

The Prisoner's Dilemma (PD) is usually the starting point of most literature devoted to study paradoxes associated with the absence of cooperation<sup>3</sup>. Nevertheless, in most real-life situations decisions are not absolutely simultaneous, but sequential. In addition, as Hirschleifer (1987, 2001) shows, most relevant economic situations are characterized by a sequential decision process in which some individuals have psychological predispositions to start relations. For Frank (1987), these predispositions are related to the emotions felt by the individual in a given context. Therefore, in situations where these predispositions are relevant, human relations take the form of a sequential commitment game. Hence, we use in our analysis a sequential PD (SPD) and not a standard PD. The SPD-version that we propose is extended to include Kropotkin's mutual aid concept. In addition, we propose an extension of the SPD to include Kant's moral rules in order to compare the implications of both moral concepts. That is, cooperation as a *choice* in Kropotkin and as a *duty* in Kant.

The plan of the paper is as follows. Section (2) presents and compares, in a non-cooperative framework, the SPD game extended to Kropotkin's and Kant's moral concepts. Section (3) analyses these games in an indirect evolution framework. We show that cooperation and generalized acceptance of Kropotkin's mutual aid principle can be an evolutionary equilibrium, while this is not possible with Kantian morality. Section (4) proofs that, in an indirect evolution framework with incomplete information, neither Kropotkin's nor Kantian moralities are likely to be the outcome since general defection will be the rule. However, section (5) shows that if some agents have an imperfect detection technology, moral agents and cooperative behavior will not disappear under Kropotkin's morality, while they will vanish with Kant's rules. Section (6) concludes.

## 2. The extended SPD-game with rational agents

The basic SPD-game is shown in figure (1) in its extensive form game. The first monetary payoff at each end node corresponds to the first mover, the second to the second mover. The two options of player  $i$  ( $i=1,2$ ) may be interpreted as "cooperation" ( $C_i$ ) and "defection" ( $D_i$ ). We assume that the payoffs are ranked as follows:  $S' < P' < R' < T'$ , and that they

---

<sup>3</sup> We will focus in our analysis on a Prisoners's Dilemma without exit option. Vanberg and Congleton (1992) have shown, by means of simulations, that a "moral program" (specified as one that never defects, but exits in response to an opponents defection) is successful in competition with a variety of alternative programs, including Tit for Tat, if an exit option is available in the Prisonier's Dilemma.

satisfy:  $(S' + T')/2 < R'$ . Without loss of generality, we normalize the payoffs so as to set  $S = 0$  (*i.e.*  $X = X' - S'$ ,  $\forall X' = S', R', T'$ ).

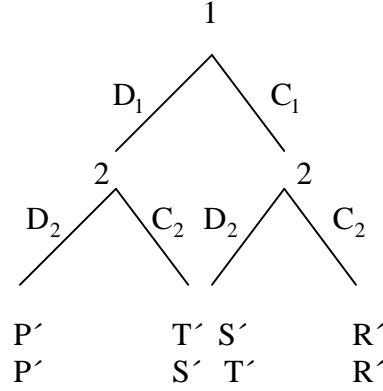


Figure (1). The SPD-game

Whatever  $j$  chooses, the best response of agent  $i$ ,  $i \neq j$ , is  $D_i$ . Hence, the only rational solution is  $(D_1, D_2)$ . This result in complete information does not depend on the timing but on the assumption that both player utilities are strictly monotonic functions of the monetary payoffs.

As stated above, Kropotkin's mutual aid concept must be able to restrain the individual "ego" (*i.e.* the egoistic character of the individual). If the agent has the mutual aid pre-disposition he will feel bad if, by privileging his selfishness, he betrays his partner while this one seeks cooperation. Thus, we will assume that the payoff  $T$  is associated with an internal feeling ( $e$ ) that will be strong for an individual with the mutual aid pre-disposition and weak (or zero) for an individual who does not have this pre-disposition. However, if the partner chooses  $D$ , Kropotkin's mutual aid idea does not imply to "take a sacrifice" for others. Kropotkin's mutual aid asks to support cooperation when this outcome is possible, but it does not mean an unconditional commitment to cooperate. Hence, the utilities are now functions of both the monetary payoff and the feeling associated with the mutual aid principle. Figure (2) shows what we will call "Kropotkin's SPD".

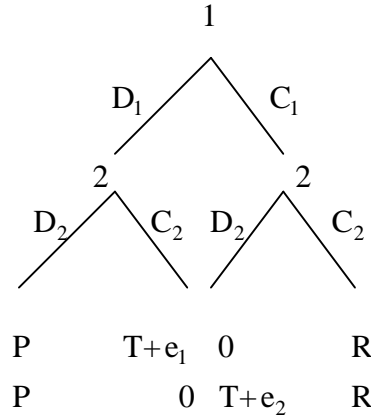


Figure (2). Kropotkin's SPD

The preference parameter  $e_i$  is internal to the actors and is not influenced by the monetary payoff. Without loss of generality, we define  $e_i$  as follows<sup>4</sup>:  $T+e_i^- < R < T+e_i^+$ ,  $i = 1,2$ . Remark that  $e_i^-$  implies a “moral” concern that changes the initial relation  $T > R$ , while  $e_i^+$  implies no moral concern, or a moral concern too small to change this initial relation between the monetary payoffs. Therefore, we will call  $e_i^-$ -agents “moral agents” (in Kropotkin's sense) and  $e_i^+$ -agents “amoral agents” (in the sense of being insufficiently motivated by morals to change their behavior).

The rational behavior in the new game depends on the values of  $e_i$ . In a situation where  $e_2 = e_2^+$ , the strategy combination  $(D_1, D_2)$  is the single subgame perfect equilibrium of the game. As usual, we start examining the game at the second stage. At this step, only agent 2 plays. Acting rationally, he chooses the strategy that ensures the highest payoff, in this case P if agent 1 plays  $D_1$ , and  $T+e_2^+$  if agent 1 chooses  $C_1$ . Consequently, agent 2's best response is to play  $D_2$ . Knowing that agent 2 will play  $D_2$ , agent 1 will choose  $D_1$  because he prefers to obtain P rather than 0. On the other hand,  $e_2 = e_2^-$  implies that the single subgame perfect equilibrium is  $(C_1, C_2)$ . In this case, players reach a Pareto-efficient result yielding a gain of R

<sup>4</sup> Our game implies, to some extent, to define the intervals most likely to support cooperation. Kirchkamp (1996), for example, developed simulations with the repeated PD where constraints are defined on the values of the profits. In this study,  $T'=1$  and  $P'=0$ , whereas the values of  $R'$  and  $S'$  vary within the limits of the intervals given by:  $0 < R' < 1$  and  $-2,5 < S' < 0$ . Kirchkamp's results are related to a great number of simulations in the  $R' \times S'$  space. The only strategies available are “always defect” and “always cooperate”. He shows that mutual cooperation appears only for  $R' > 0,7$  and  $S' > -1$ , while universal defection prevails otherwise. Simulations by other authors confirm the result that cooperation can be promoted by suitable modifications of the profits. To obtain this result, the values of  $R'$  and  $S'$  must be as high as possible within the constraints that define the PD.

for each. Thus, we obtain a strategy similar to Tit For Tat (Axelrod, 1980), *i.e.* agents answer cooperation with cooperation and defection with defection.

We will now exceed the requirements of the mutual-aid principle and evoke a more demanding framework: Kant's (1785/1965) categorical imperative. According to Kant, it is possible to determine moral principles able to control our behavior towards others. Each agent must put himself at the place of the others so as to identify his duties towards them and, *a fortiori*, towards himself. Whatever his objective may be, whatever desires he may feel, whatever inclinations he may have, if he wants to act morally, he has to conform his act to the categorical imperative (or the "imperative of duty"): "Act only on that maxim whereby thou canst at the same time will that it should become a universal law" (Kant, 1785/1965). Since the outcome of generalized cooperation is preferred to the one associated to generalized defection, players would wish cooperation to become a universal law (Sugden, 1991). Thus, they have to cooperate. Therefore, cooperation is a duty and not a choice in response to a given environmental situation. In the SPD that means, for both players, that "bad feelings" do not only appear if they betray a partner who cooperates, but in all situations where they do not follow their moral duty to cooperate. Thus, the parameter ( $e$ ) also modifies the utility  $P$  that agents obtain if both choose defection. Figure (3) gives the extensive form of what we will call a "Kantian SPD".

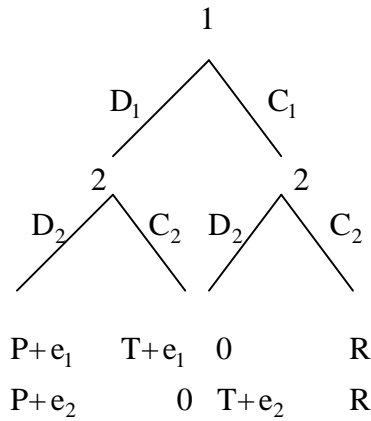


Figure (3). Kantian SPD

The assumptions of the game are now:  $T + e_i^- < R < T + e_i^+$  and  $P + e_i^- < 0 < P + e_i^+$ ,  $i = 1, 2$ . In this context,  $e_i^-$  implies a moral concern that changes the initial relation  $T > R$  and also the relation  $P > 0$ , while  $e_i^+$  implies no moral concern, or a moral concern too small to change this



initial relation between the payoffs. Thus, we will call an  $e_i^-$ -agent a “moral” agent (in Kant’s sense) and an  $e_i^+$ -agent an “amoral” agent (as in the previous game).

It is easy to show that if  $e_2 = e_2^+$ , then: (i) for  $e_1 = e_1^+$ ,  $(D_1, D_2)$  is the only subgame perfect equilibrium; (ii) for  $e_1 = e_1^-$ ,  $(C_1, D_2)$  is the only subgame perfect equilibrium. On the contrary, if  $e_2 = e_2^-$ , then: (i) for  $e_1 = e_1^+$ ,  $(D_1, C_2)$  is the only subgame perfect equilibrium; (ii) for  $e_1 = e_1^-$ ,  $(C_1, C_2)$  is the only subgame perfect equilibrium. As it can be seen, the actions of the two types of agents are only determined by their own type, and have no relation with the actions undertaken by their partner. Thus, we can write the following remark:

### **Remark 1**

*In a Kantian SPD, a moral agent ( $e_i^-$ -agent) will always cooperate and an amoral agent ( $e_i^+$ -agent) will always defect.*

That is, a “moral” agent in the sense of Kant (an  $e_i^-$ -agent) will always do what he would like to be the universal law, *i.e.* to cooperate since this would yield the highest benefit if followed by all players. However, as a direct result of remark 1, mutual cooperation only appears if both agents are  $e_i^-$ . Therefore, with Kant’s categorical imperative all the outcomes of the game can be equilibrium solutions, depending upon the types of the agents involved in the game.

Within the non-cooperative framework we could, modifying the game, ensure cooperation irrespective of players' characteristics. We could, for instance, assume the existence of norms (Bendor and Swistak, 2002) or other forms of external punishment. However, if we include norms in the analysis we exclude *de facto* the possibility that agents develop a cooperative preference due to internal moral concerns. However, experimental economics has shown that such a preference for cooperation may in fact exist (Charness and Rabin, 2002; Sandbu, 2002). Thus, we should only take into account elements which are internal to the agents. Our aim is to explain if Kropotkin’s and/or Kant’s moral rules can favor global cooperation, and under which conditions, and the evolutionary game-theory framework is well adapted for this task, especially in a context with indirect evolution. This latter approach, applying the ESS concept to preferences rather than to strategies, permits to determine preferences endogenously.

### 3. Indirect evolution

Since we are now in an evolutionary game-theory framework we assume a large random-mixing population. Agents engage only in sequential symmetric pairwise contests. One player is labeled player 1 and the other player 2.

Define  $E_i = \{e_i^-, e_i^+\}$  as the set of values that the parameter  $e_i$  of agent  $i$  can take ( $i = 1, 2$ ). Each one of these two values describes a type of agent. That is, as in the previous section, agents are either amoral ( $e_i^+$ ) or moral ( $e_i^-$ ), in Kropotkin's or Kant's sense (depending on the game under consideration). To ensure a symmetric game, the position in the game of both players is random. Thus, Nature (agent 0) defines this position. Without loss of generality, we consider that the probability for each agent to play first is  $\frac{1}{2}$ . The symmetrical extensive form associated with this new version of the game is described by figure (4) for Kropotkin's SPD. We will use the ESS concept, as defined by Maynard Smith and Price (1973), since it is one of the most widely used concepts of evolutionary game theory.

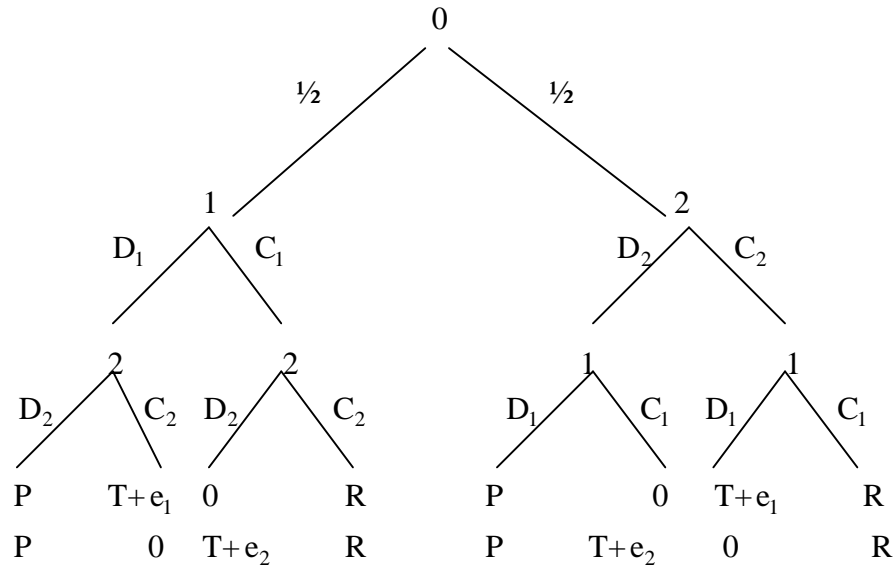


Figure (4). Symmetric Kropotkin's SPD

#### Definition 1

Given two strategies  $e_1$  and  $e_2$ ,  $U(e_1, e_2)$  is the utility that agent 1 obtains by choosing  $e_1$  while the other player chooses  $e_2$ . A monomorphic population (i.e, when all components of

the population use the same strategy) is stable if and only if all agents adopt a strategy  $e_1$  that meets the two following conditions (for  $i = 1, 2$ ):

1. *Equilibrium condition:*  $U(e_1, e_1) \geq U(e_2, e_1) \quad \forall e_2$
2. *Stability condition:* if  $e_2 \neq e_1$  and  $U(e_1, e_1) = U(e_2, e_1)$ , then  $U(e_1, e_2) > U(e_2, e_2)$ .

## Definition 2

If  $e_1$  satisfies conditions (1) and (2), we say that  $e_1$  constitutes an ESS.

That is, a given behavior is evolutionary stable within a population if no other strategy can make a better score playing against  $e_1$  than  $e_1$  itself, and, at the same time, if  $e_1$  cannot do better against itself than  $e_2$ , then  $e_1$  makes a better score against  $e_2$  than  $e_2$  itself.

Thus, the analysis of the game consists in studying the evolutionary stability of the preference parameter "e". Consequently, the evolutionary approach consists in determining  $U(e_1, e_2)$  for  $e_1 \in \{e_1^-, e_1^+\}$  for agent 1 and  $e_2 \in \{e_2^-, e_2^+\}$  for agent 2. Note that this discussion is analogous to section (2). However, results depend not only on each agent's type, but also on its position within the game. Thus, we obtain the following bimatrix, showing the expected gain of each type of agent.

|         |                   |                   |
|---------|-------------------|-------------------|
|         | $e_2^-$           | $e_2^+$           |
| $e_1^-$ | R,R               | (P+R)/2 , (P+R)/2 |
| $e_1^+$ | (P+R)/2 , (P+R)/2 | P,P               |

Figure (5) Expected gains in Kropotkin's symmetric SPD

This bimatrix is obtained repeating the following reasoning for the different combinations of "e". Consider the case where  $e_1 = e_1^-$  and  $e_2 = e_2^-$ . With a probability  $1/2$  agent 1 will open the game. In this case, the perfect equilibrium will be  $(C_1, C_2)$ , and gains (R,R). If agent 2 opens the game, the perfect equilibrium will also be  $(C_1, C_2)$ , and gains also (R,R). Thus  $U(e_1^-, e_2^-) = R/2 + R/2 = R$ . Since  $R > P$ , it is easy to show that the payoffs of 1 are higher for  $e_i = e_i^-$  than for  $e_i = e_i^+$ . Hence, we can write:

**Proposition 1**

*When players can identify their opponent's type in Kropotkin's symmetric SPD,  $e_i^-$  is the only evolutionarily stable strategy.*

*Proof: direct from figure (5).*

According to proposition (1), starting from a situation where moral (in Kropotkin's sense) and amoral agents exist, moral agents will impose themselves and morality will spread. The drawback is that this result is related to the fact that information is complete. The first mover knows the characteristics of his partner, and will thus adapt his strategy. Since encounters between moral agents yield a higher utility than encounters between amoral agents, the learning process will imply for amoral agents to develop morality.

This result is not completely new. Within a different framework, Gauthier (1986) arrived to a similar conclusion. One of Gauthier's interesting intuitions is to suppose that it can be rational for agents to choose certain pre-dispositions before the game starts. These pre-dispositions act as constraints on the set of possible actions of the actors. Gauthier considers two types of strategies: individual strategies associated to players maximizing their individual profit, and joint strategies associated to players maximizing the collective profit. The model also includes two types of actors: (i) "direct" maximizers, who only seek to maximize their own utility (these agents betray in the one-shot PD), and (ii) "constraint" maximizers (they choose a joined strategy if the profit is equal to or higher than the one associated to an individual strategy). Introducing these two behaviors in a static PD he shows that choosing a constraint maximization is a rational act, which supports mutual cooperation. Our result is equivalent, although we consider a SPD and do not consider pre-dispositions as given. However, the similarity of the results obtained suggests that cooperation is not due to sequentiality but to complete information.

It is interesting to remark that the evolutionary framework analyzed in this paper does yield a similar, but not identical, result to Becker's (1974) "rotten kid theorem". Becker studies a game with two players where one of the actors is altruistic. He shows that even if the recipient of the altruistic behavior adopts an egoistic behavior, he will not harm his benefactor. Suppose that an action increases the income of the recipient and reduces that of the altruist. The altruist, taking into account the reduction of his income, may reduce his contribution more than what the recipient would gain. In other words, even an egoist

internalizes the benefit that he can obtain from his benefactor. That is, the altruist pushes indirectly the other individual to maximize the altruist's income, but the latter does not become an altruist (Becker (1974) calls this the "rotten kid theorem"). This result is based on given preferences and individuals, taking into account their preferences, can only choose the strategies imposed by their rationality. In the indirect evolution framework, individuals have a broader choice, they can build their preferences. By learning from others, they can determine the best preferences to establish a social link based on reciprocity and cooperation. Evolutionary forces will determine not only the behavior but also the non-monetary preferences that are more convenient. With complete information the answer is clear: morality in Kropotkin's sense (*i.e.* mutual-aid) is better than amorality (*i.e.* selfishness). Thus, by permitting egoists to re-examine their concept of social relations, they will (progressively) become more and more moral (in Kropotkin's sense).

The symmetrical extensive form of the SPD with agents guided by Kant's categorical imperative is shown in figure (6). Following the same type of reasoning as for figure 5, we obtain figure (7). However, using remark 1 the same result can be obtained more directly.

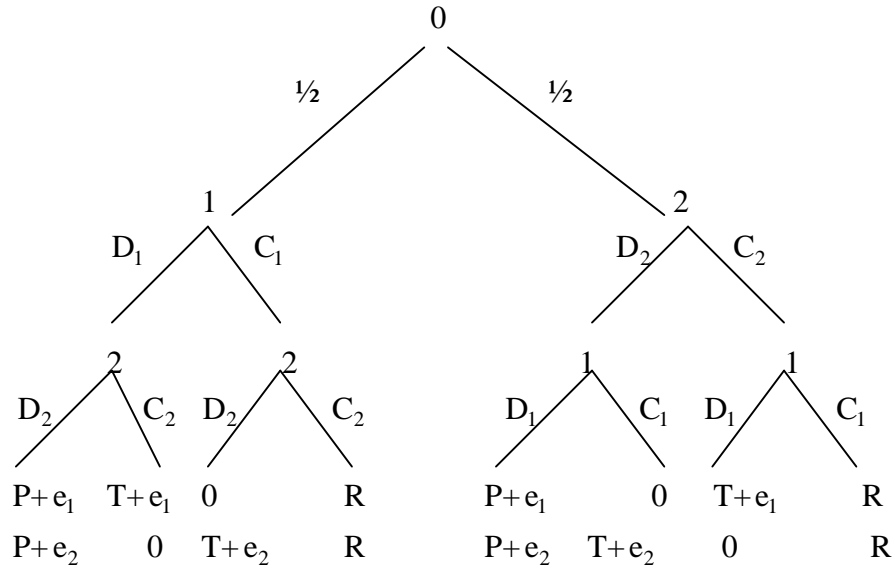


Figure (6) Symmetric Kantian SPD

|         |         |         |
|---------|---------|---------|
|         | $e_2^-$ | $e_2^+$ |
| $e_1^-$ | R,R     | 0,T     |
| $e_1^+$ | T,0     | P,P     |

Figure (7) Expected gains in the Kantian symmetric SPD

This permits us to write the following proposition:

### **Proposition 2**

*When players can identify their opponent's type in a symmetric Kantian SPD,  $e_i^+$  is the only evolutionarily stable strategy.*

*Proof: direct from figure (7).*

That is, although moral agents in the Kantian SPD have a higher level of commitment in favor of cooperation than in Kropotkin's SPD, the Kantian SPD brings us back to the standard prisoner's dilemma. Hence, the Kantian requirement, by setting cooperation as a duty, encourages others to defection and not to cooperation and obtains the opposite result to that what it pretended to promote (*i.e.* in equilibrium all agents are amoral ( $e_i^+$ )). In order to study the robustness of our results, we are now going to analyze the same models with incomplete information. This will enable us to find a strategic dimension in the sequentiality.

### **4. Indirect evolution with incomplete information about the partner's type**

Let  $q$  denote the share of  $e^-$  type agents in the population and  $(1-q)$  the share of  $e^+$  type agents, with  $0 \leq q \leq 1$ . We assume that the probability distribution  $(q, 1-q)$  is common knowledge. Moreover, we assume that players know their own type, but not the type of their partners. The probability of playing first (or second) is for each player  $\frac{1}{2}$  (independently of their type).

In Kropotkin's SPD, the optimal decision of the second mover does not depend on the type of the first, but on his behavior. Thus, the second mover's information set becomes strategically irrelevant and we can predict the second player's behaviour. If he is an  $e_i^-$  type agent (a moral agent in Kropotkin's sense), he will use strategy  $C_i$  if the first player chooses  $C_j$ , and strategy  $D_i$  if the first player's choice is  $D_j$ ,  $i \neq j$ . If he is a  $e_i^+$  type agent (an amoral agent), he will always choose  $D_i$ ,  $i = 1, 2$ .

By retracing the development of the game one stage we examine the choice made by the first mover. Irrespective of his type, the first mover chooses to cooperate if:  $q > P/R$ . Suppose that an agent 1 of type  $e_1^-$  chooses  $C$ . There is a probability  $q$  that agent 2 will be of

type  $e_2^-$ . Since a  $e_2^-$ -agent 2 will adopt strategy  $C_2$ , agent 1 will obtain R with a probability q. If faced with an agent 2 of type  $e_2^+$  (with probability 1-q), agent 1 will receive 0. If agent 1 chooses  $D_1$ , the gain would be P if faced with an agent 2 of type  $e_2^-$  (with probability q), but also when matched with an agent 2 of type  $e_2^+$  (with probability 1-q). Therefore,  $C_1$  is more beneficial than  $D_1$  if  $qR > P$ . Since the same results can be found for an  $e_1^+$  agent playing first (owing to the fact that T does not intervene in the argument) the first mover will choose  $C_i$  if  $q > P/R$  regardless of his type. Conversely, he will prefer  $D_i$  if:  $q < P/R$ .

However, if we wish to characterize the expected utility by each of these two types, we must study if the parameter values are linked to the first or the second condition as determined above. Consider  $q > P/R$ . With probability  $1/2$  an agent  $i$  of type  $e_i^-$  ( $i = 1,2$ ) will be assigned the role of first mover. By choosing  $C_i$  (owing to the initial condition) his gain is  $qR$ . However, with probability  $1/2$  he will play the game as a second mover. We know that an agent  $i$  of type  $e_i^-$  playing second will choose  $C_i$  rather than  $D_i$  (obtaining R). Thus, the expected gain for an  $e_i^-$  agent is:  $1/2[qR + R]$ . If agent  $i$  is of type  $e_i^+$ , he will play first with probability  $1/2$  and get a gain of  $qR$ . If he plays second (with probability  $1/2$ ), it will be more beneficial for him to choose  $D_i$  rather than  $C_i$ , since this choice yields a gain of T. Hence, the expected gain for type  $e_i^+$  is:  $1/2[qR + T]$ . Thus, the expected gain associated with type  $e^+$  is superior to the one of type  $e^-$  (given  $T > R$ ). If we adopt a similar reasoning for the case where  $q < P/R$ , we note that the first mover will choose D, and so will do the second mover (and both would obtain a benefit equal to P). Therefore, we can write:

### **Proposition 3**

*In Kropotkin's SPD with incomplete information concerning agents' type, no population containing a proportion  $1 > q > P/R$  of  $e^-$  type agents can be evolutionarily stable.*

According to proposition (3), if amoral agents appear in a population formed only by moral agents (in Kropotkin's sense) the cooperative behavior will be significantly reduced. As amoral agents appear, they will be more and more encouraged since they obtain higher (monetary) payoffs than moral agents. However, if the proportion of moral agents is lower than  $P/R$ , the evolutionary pressure plays no role anymore and first movers will always

choose D. Thus, moral agents (with the mutual aid pre-disposition) will not disappear; they will just become a minority. Nevertheless, the outcome will be generalized defection (since even moral agents will choose to betray in this environment dominated by amoral agents).

In the case where  $q < P/R$  we are confronted with a problem (since payoffs are equal). However, Selten (1988) proposes to modify the game slightly in such a case. The idea is that, in a complex universe, agents may make errors in choosing their behavior. This implies to substitute the concept of ESS with that of limit ESS (LESS), which associates errors with probabilistic disturbances.

Consider a modified game where strategy C can be adopted with a probability  $\alpha$  by the player who opens the game, where  $\alpha \rightarrow 0$ . The expected gain for a  $e^-$ -agent is now  $E(e_1^-) = (1-\alpha)P + \alpha(qR + R)/2$ . For an  $e^+$ -agent the expected gain is  $E(e_1^+) = (1-\alpha)P + \alpha(qR + T)/2$ . Since, by assumption,  $R < T$ , the expected profit of a  $e^+$ -agent is higher than the one of a  $e^-$ -agent. Thus, as long as the strategy C can be selected (even with an arbitrarily small probability) only a  $e^+$ -monomorphic population (only amoral agents) is evolutionary stable, within the LESS concept. Hence:

#### **Proposition 4**

*In Kropotkin's SPD with incomplete information concerning agents' type, only a  $e^+$ -monomorphic population is limit evolutionarily stable.*

Thus, information is the main element determining agents' behavior in Kropotkin's SPD. If the individuals know each others type before the game starts, then cooperation is evolutionary stable. On the other hand, with incomplete information or, more precisely, when agents do not know precisely the type of their partners, the first mover will not take the risk to cooperate. That is, evolution supports selfishness if information is not complete. In this context, selfishness is not the will to benefit from the benevolence of others; it is more a way of protection against the uncertainty related to asymmetric information. Thus, selfishness is more a form of fear.

Retracing the same argument for the Kantian SPD we will start by analyzing the decision of the first mover. In the case of the Kantian SPD, if an  $e_1^-$  playing first chooses  $C_1$  he will obtain  $qR$ , while choosing  $D_1$  he will get  $[qT + (1-q)P]$ . Thus, an  $e_1^-$  agent will always choose  $C_1$  irrespective of the level of  $q$  (recall that an  $e_1^-$  agent prefers  $R$  over  $T$ ), while a



$e_1^+$  agent will choose  $D_1$  irrespective of the level of  $q$ . That is, we find the same result as for complete information. In other words, the asymmetry of information plays no role in this game. If we now compare the expected payoffs of a  $e_1^-$ -agent and an  $e_1^+$ -agent, knowing that they have both a probability  $1/2$  of playing first, we get  $E(e_1^-) = qR$  and  $E(e_1^+) = qT + (1-q)P$ . Since  $E(e_1^-) < E(e_1^+)$  regardless of  $q$  we can write:

### **Proposition 5**

*In a Kantian SPD with incomplete information concerning agents' type, no population containing a proportion  $1 > q > 0$  of  $e^-$  type agents can be evolutionary stable.*

That is, as long as moral agents (in Kant's sense) and amoral agents coexist the population will not be evolutionary stable. Kantian moral agents will systematically loose compared to amoral agents and Kantian morality will be progressively abandoned. The reason is that, as in complete information, moral agents ( $e^-$  agents) cooperate systematically (recall remark 1), while amoral ( $e^+$  agents) defect systematically (and this implies a higher monetary payoff for amoral agents). Thus, the only situation where Kantian behavior can lead to global cooperation is when all agents are moral ( $e^-$  agents), the situation assumed by Sugden (1991) and promoted by Kant himself.

## **5. Indirect evolution with incomplete information about the partner's type and imperfect type detection**

In the previous sections we have shown that, within Kropotkin's SPD, generalized cooperation will be the outcome if agents are able to identify their opponent's type, while cooperative behavior will disappear if they are unable to do so. We will now analyze the outcome if some agents are able to identify their opponent (informed players  $I$ ), while others are not capable to find out the type of their partners (uninformed players  $U$ ). The share of informed players will be denoted by  $l$  and that of uninformed players by  $(1-l)$ . However, since finding out the moral type of the opponent is a difficult task we assume that informed players are only able to obtain an imperfect signal about the type of their counterpart. We assume that if the agent playing second is moral ( $e^-$ ), the imperfectly informed player playing first (in the case of the second player this information is irrelevant) gets a signal indicating that his partner

is moral with probability  $\mu^-$  and a signal indicating that he is amoral with probability  $(1 - \mu^-)$ . If the second player is  $e^+$ , the informed player gets a signal indicating an amoral agent with probability  $\mu^+$  and a signal indicating a moral agent with probability  $(1 - \mu^+)$ . In order to have a meaningful detection technology we require  $1/2 < \mu^-, \mu^+ \leq 1$ . To obtain this information the informed player has a fix cost  $F \geq 0$ , which we assume lower than the level  $F_0$  defined in the Appendix (above which no agents would use the detection technology). Informed players use this information to update their prior beliefs using Bayes's rule.

This 'detection technology' is similar to the one proposed in Güth and Kliemt (1998) and Güth *et al.* (2000) and includes all kind of means to obtain information, such as the physical aspect of the opponent, his or her look, asking a given set of questions, hiring detective services, etc. In fact, most emerging social groups set up a given set of signs to facilitate the detection of other agents with the same type of religion, social preference and, why not, moral predispositions. This practice has been the case with more or less secret societies (e.g. the masons), which imply a set of moral conducts, but also with religions (e.g. the first Christians in the times of Rome), which do also imply a moral code. That is, the type of detection technology that we are postulating can be seen as an imperfect way to detect people of a given moral code, maybe asking them a few questions which are supposed to be know only by this particular type of morality. Of course, and as is obvious from the examples just given, these kind of 'technologies' are imperfect by their own nature.

What we have now is actually a system where agents can be moral or amoral and, at the same time, informed or uninformed. Both characteristics evolve in an indirect evolution framework. This system can be shown (see Appendix) to follow the phase-diagram in figure (8) or the phase diagram in figure (9) depending upon the relationship between  $\mu^-(R - P)$  and  $(1 - \mu^+)(T - P)$ . If  $\mu^-(R - P) > (1 - \mu^+)(T - P)$  the system follows figure (8) and if this relation is reversed the system follows figure (9). That is, we have to compare (i) the probability of correctly detecting a moral agent  $\mu^-$  times the net reward of mutual cooperation over mutual defection  $(R - P)$  with (ii) the probability of getting a sign of morality from an amoral agent  $(1 - \mu^+)$  times the net benefit of treating as compared to mutual defection  $(T - P)$ . That is, the better our detecting technology is, the closer we are from figure (8) and, the higher the reward of treason as compared to mutual cooperation is, the closer we are from figure (9).

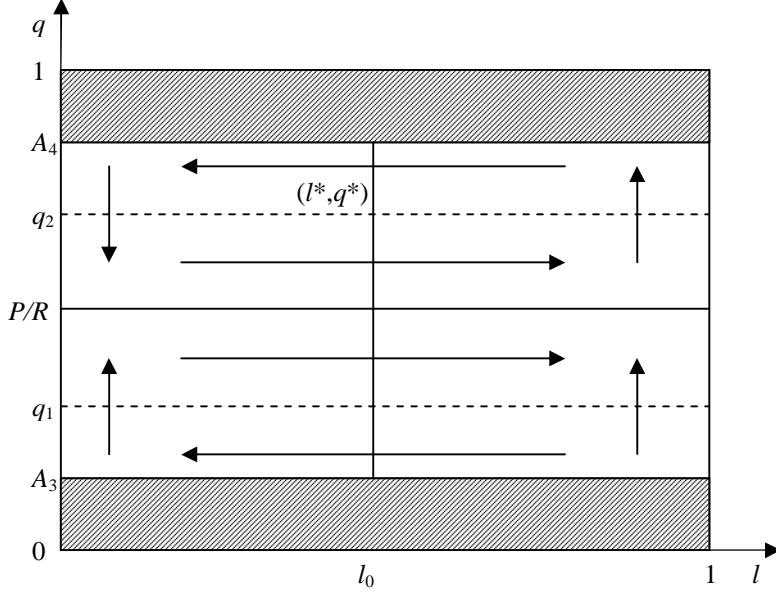


Figure (8) First phase-diagram for Kropotkin's SPD with imperfect type detection

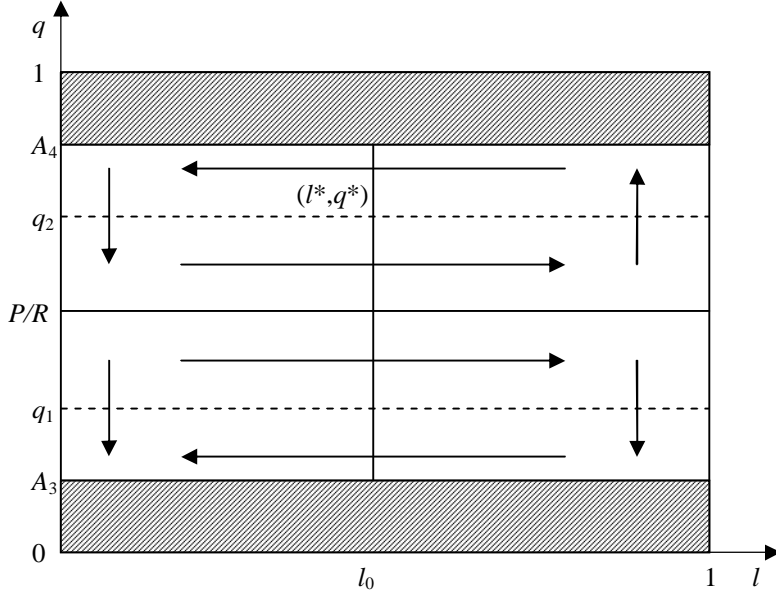


Figure (9) Second phase-diagram for Kropotkin's SPD with imperfect type detection

The dynamics described in figures (8) and (9) can be summarized as follows (the expressions for the limits  $A_3$  and  $A_4$  can be found in the Appendix; below, respectively above, these limits informed and uninformed agents behave in the same manner):

- i. Above the limit  $A_4$  and under  $A_3$  agents behave as uninformed agents (see section 4). If the initial point is above  $A_4$  this will eventually bring the population to the region discussed next, since the number of moral agents will tend to decrease. If the initial population of moral agents is below  $A_3$  general defection will be the outcome (see section 4).

- ii. In the region below  $A_4$  and above  $P/R$ , the system will perpetually cycle around the degenerated rest point  $(l^*, q^*)$ , see Appendix. That is, the number of moral agents will constantly change, but they will not disappear. In this region uninformed moral agents will cooperate (see section 4) and informed moral agents will also cooperate if they receive a signal  $e^-$  (see Appendix). Thus, cooperative behavior will not disappear.
- iii. In the region below  $P/R$  and above  $A_3$ , the behavior of the system depends upon the relationship between  $\mu^-(R-P)$  and  $(1-\mu^+)(T-P)$ . If  $\mu^-(R-P) > (1-\mu^+)(T-P)$ , figure (8), starting in this region moral agents will tend to increase. Informed agents will cooperate in this region if they receive a  $e^-$  signal (see Appendix), while uninformed agents will not cooperate (see section 4). However, the increase in moral agents will eventually allow the system to reach the region discussed in (ii). If  $\mu^-(R-P) < (1-\mu^+)(T-P)$ , figure (9), moral agents will tend to disappear and the lower region discussed in (i) will eventually be reached.

In the previous section we have shown that in Kropotkin's SPD with incomplete information moral agents will not disappear, but that they will become such a minority that even moral agents will choose to betray, so that cooperative behavior will disappear. Including an imperfect detection technology we have shown the conditions under which, not only moral agents do not disappear, but cooperative behavior does also not disappear. Thus:

**Proposition 6**

*In Kropotkin's SPD with incomplete information and imperfect type detection, and starting from a significant number of moral agents [  $q > A_3$  if  $\mu^-(R-P) > (1-\mu^+)(T-P)$  and  $q > P/R$  if  $\mu^-(R-P) < (1-\mu^+)(T-P)$  ], moral agents ( $e^-$ ) will not disappear and cooperative behavior will also not disappear.*

*Proof: see Appendix and figures (8) and (9).*

Nevertheless, as the streamlines in figures (8) and (9) show the equilibrium of the system is a vortex. Thus, we have an unstable equilibrium and no proportion of moral agents will be evolutionary stable (except if the initial population of moral and informed agents is precisely the vortex).

Analyzing the impact of an imperfect detecting technology in the Kantian SPD has no interest, since cooperation is not achieved neither under complete nor under incomplete information. That is, even with an imperfect detection technology moral agents will tend to disappear and cooperative behavior will disappear with them.

## **6. Conclusion**

With complete information, Kropotkin's mutual-aid concept implies a strategic behavior whose main features are similar to the Tit For Tat strategy. Cooperation is a rational choice and cooperation is answered with cooperation, yielding an evolutionary stable situation where all agents follow Kropotkin's moral and cooperate. However, applying Kant's moral rules, and assuming cooperation as a duty, generalized defection is the outcome. In other words, morality can be the solution to the Sequential Prisoner's Dilemma as long as the cooperative behavior is seen as a strategic choice (Kropotkin) and not as an intangible duty imposed to individuals (Kant).

When information is not complete, selfish behavior is the best way of protecting oneself against non-cooperative behaviors and cooperation will not be achieved with any of the moral concepts discussed in this paper. However, if some agents are able to detect the moral type of their partners, and even if this detection is not perfect, moral agents in Kropotkin's sense will not disappear if they are sufficiently important in the initial population, and cooperative behavior will also not disappear. This does not apply to moral agents in Kant's sense, since they will tend to disappear in all the environments considered, and cooperative behavior will vanish with them.

## **References**

- Axelrod, R. (1980), "Effective Choice in the Prisoner's Dilemma", *Journal of Conflict Resolution*, 24, p. 3-25.
- Ballet, J. (2000), "Altruisme et Biens Collectifs, une Présentation de la Littérature", *Revue Economique*, 4, p. 789-811.
- Becker, G.S. (1974), "A Theory of Social Interactions", *Journal of Political Economy*, 82, p. 1063-1093.
- Becker, G.S. (1981), *A Treatise on the Family*, Cambridge Mass., Harvard University Press.

- Bendor, J., Swistak, P. (2002), "The Evolution of Norms", *American Journal of Sociology*, 106, p. 1493-1545.
- Bergstrom, T. C. (1995), "On the Evolution of Altruistic Ethical Rules for Sibling", *American Economic Review*, 85, p. 789-811.
- Berninghaus, S., Güth, W., Kliemt, H. (2003), "From Theology to Evolution: Bridging the Gap between Rationality and Adaptation in Social Explanation", *Journal of Evolutionary Economics*, 13, p. 385-410.
- Bordignon, M. (1990), "Was Kant Right? Voluntary Provision of Public Goods under the Principle of Unconditional Commitment", *Economic Notes*, 3, P. 342-372.
- Charness, G., Rabin, M. (2002), "Understanding Social Preferences with Simple Tests", *Quarterly Journal of Economics*, 114, p. 817-868.
- Etzioni, A. (1987), "Toward a Kantian Socio-Economics", *Review of Social Economy*, 45, p. 37-47.
- Dekel, E., Ely, J.C., Yilankaya, O. (2007). "Evolution of Preferences," *Review of Economic Studies*, 74(3), p. 685-704.
- Fong, C. M., Bowles, S., Gintis, H. (2006), "Strong Reciprocity and the Welfare State." In S.-C. Kolm and J. M. Ythier (ed.) *Handbook on the Economics of Giving, Altruism and Reciprocity*, Volume 2, Amsterdam, Elsevier, p. 1440-64
- Frank, R.H. (1987), "If Homo Economicus Could Choose His Own Utility Function, Would He Choose One with a Conscience?", *The American Economic Review*, 77(4), p. 593-604.
- Gauthier, D. P. (1986), *Morals by Agreement*, Oxford, Oxford University Press.
- Güth, W., Kliemt, H. (1994), "Competition or Co-Operation: on the Evolutionary Economics of Trust, Exploitation and Moral Attitudes", *Metroeconomica*, 45, p. 155-187.
- Güth, W., Kliemt, H. (1998), "The indirect Evolutionary Approach: Bridging the Gap Between Rationality and Adaptation", *Rationality and Society* 10(3), p. 377-399.
- Güth W., Kliemt H., Peleg, B. (2000), "Co-evolution of Preferences and Information in Simple Games of Trust", *German Economic Review*, 1(1), p. 83-110.
- Güth, W., Napel, S. (2006). Inequality aversion in a variety of games: an Indirect Evolution Analysis". *The Economic Journal*, 116, 1037–1056.
- Güth, W., Yaari, M. (1992), "An Evolutionary Approach to Explaining Reciprocal Behavior in a Simple Strategic Game", in U. Witt (ed.), *Explaining Process and Change – Approaches to Evolutionary Economics*, Ann Arbor, The University of Michigan Press, p. 23-34

- Guttman, J. M. (2000), "On the Evolutionary Stability of Preferences for Reciprocity", *European Journal of Political Economy*, 16, p. 31-50.
- Hirshleifer, J. (1987), "On the Emotions as Guarantors of Threats and Promises", in J. Dupre (ed.), *The Latest on the Best*, Cambridge, Mass., MIT Press, p. 307-326
- Hirschleifer, J. (2001), "Game-Theoretic Interpretations of Commitment", in R.M. Neese, *Natural Selection and the Capacity for Commitment*, New York, Russell Sage Press, p. 77-92.
- Kant, E. (1785/1965), *Grundlegung zur Metaphysik der Sitten*. Verlag von Felix Meiner, Hamburg.
- Kirchkamp, O. (1996), "Spatial Evolution of Automata in the Prisoners' Dilemma", in K.G. Troitzsch, U. Mueller, G.N. Gilbert and J.E. Doran (eds.), *Social Science Microsimulation*, Berlin, Springer-Verlag, p. 307-358.
- Königstein, M., Müller, W. (2000), "Combining Rational Choice and Evolutionary Dynamics: The Indirect Evolutionary Approach", *Metroeconomica*, 51, p. 235-256.
- Kropotkin, P. (1902), *Mutual Aid: A Factor of Evolution*, New York, Double Day.
- Laffont, J.-J. (1975), "Macroeconomics Constraints, Economic Efficiency and Ethics: An Introduction to Kantian Economics", *Economica*, 42, p. 430-437.
- Margolis, H. (1982), *Selfishness, Altruism and Rationality*, Cambridge, Cambridge University Press.
- Maynard Smith, J., Price, G. R. (1973), "The Logic of Animal Conflicts", *Nature*, n° 246, p. 15-18.
- Ok, E.A., Vega-Redondo, F., On the Evolution of Individualistic Preferences: An Incomplete Information Scenario. *Journal of Economic Theory* 97, 231\_254
- Roemer, J.E. (1996), *Theories of Justice*, Cambridge Mass., Harvard University Press.
- Sandbu, M.E. (2002), "The Road Not Taken: A Theory of Set-Dependent Fairness Preferences", *Working Paper*, Harvard University.
- Selten, R. (1988), "Evolutionary Stability in Extensive Two-Person Games: Correction and further Development", *Mathematical Social Sciences*, 16, p. 223-266.
- Sugden, R. (1991), "Rational Choice: A Survey of Contributions from Economics and Philosophy", *Economic Journal*, 101, p. 751-758.
- Tuomela, R. (1990), "What are Goals and Joint Goals", *Theory and Decision*, 28, p. 1-20.
- Tuomela, R., Miller, K. (1988), "We-intentions", *Philosophical Studies*, 53, p. 367-389.
- Vanberg, V.J., Congleton, R.D. (1992), Rationality, Morality and Exit", *American Political Science Review* 86(2) 418-431.

# Appendix<sup>1</sup>

In Kropotkin's SPD with incomplete information and imperfect detection type, the probability of receiving a  $e^-$ -signal is  $[\mu^- q + (1 - q)(1 - \mu^+)]$  and the *a posteriori* probability of meeting a  $e^-$ -type is :

$$\frac{\mu^- q}{\mu^- q + (1 - q)(1 - \mu^+)} = A_1$$

The probability of receiving a  $e^+$  signal is  $[(1 - \mu^-) q + (1 - q)\mu^+]$  and the *a posteriori* probability of meeting a  $e^+$  type, after receiving a signal  $e^+$ , is :

$$\frac{(1 - \mu^-) q}{(1 - \mu^-) q + (1 - q)\mu^+} = A_2$$

Suppose that the informed player ( $I$ ) gets a signal  $e^-$ . If he cooperates the probability of being actually faced with a  $e^-$  agent is  $A_1$  (expected gain :  $A_1 R$ ). If he does not cooperate he knows that player 2 will always betray so that the gain will be  $P$ . Thus, he will cooperate if :

$$A_1 > \frac{P}{R} \iff q > \frac{(1 - \mu^+) P}{\mu^- (R - P) + (1 - \mu^+) P} = A_3$$

Suppose now that player  $I$  gets a signal  $e^+$ . Defection will be chosen if :

$$A_2 < \frac{P}{R} \iff q < \frac{\mu^+ P}{(1 - \mu^-)(R - P) + \mu^+ P} = A_4$$

Remark that with  $1 > q > 0$  and  $1/2 < \mu^-, \mu^+ \leq 1$  we have  $1 > A_1 > A_2 > 0$  or  $0 < A_3 < A_4 < 1$ .

It is easy to show that the differentiation between informed ( $I$ ) and uninformed ( $U$ ) agents can only take place in the region where  $\frac{P}{R} \in (A_2, A_1) \iff q \in (A_3, A_4)$ . To see this, assume that an  $I$  player ignores the signal  $e^-$  and does not cooperate since  $A_1 < P/R$ . Given our assumptions this implies  $P/R > q$ , so that an  $U$  player would also not cooperate. Moreover, since  $A_1 > A_2$  an  $I$ -player should follow the  $e^+$  signal and choose defection, like an  $U$ -agent would do. The same reasoning can be applied to the case where  $A_2 > P/R$ . Thus, we focus on the region  $q \in (A_3, A_4)$ .

The dynamics of  $l$  is calculated comparing the success of  $I$ -agents as compared to  $U$ -agents. For  $q$ , the success of  $e^-$  and  $e^+$ -agents is compared. We will start with the dynamics of the informed players ( $l$ ).

The expected gain of an  $e_1$ - $I$ -type (+ or -) playing first (with probability 1/2) is  $EG(e_1-I) = \frac{1}{2} \{q [\mu^- R + (1 - \mu^-) P] + (1 - q) \mu^+ P\} - F$ . The expected

---

<sup>1</sup>The resolution method is similar to the one used in Güth *et al.* (2000).



gain of a  $e_1$ - $U$ -type (+ or -) playing first (with probability 1/2) depends on the level of  $q$ . For  $q > P/R$  the first player will cooperate and get  $EG(e_1-U) = \frac{1}{2}qR$ . If  $q < P/R$  he will defect and get  $EG(e_1-U) = \frac{1}{2}P$ . The benefit of a  $e_1^+$ - $U$ -type and a  $e_1^+$ - $I$ -type (or a  $e_1^-$ - $U$ -type and a  $e_1^-$ - $I$ -type) is the same, since the benefit depends only on the moral type, and not on the fact of being informed or uninformed.

For  $q < P/R$  we have :

$$EG_I(q) > EG_U(q) \iff q > \frac{(1 - \mu^+) P + 2F}{(1 - \mu^+) P + \mu^- (R - P)} = q_1$$

and for  $q > P/R$  :

$$EG_I(q) > EG_U(q) \iff q < \frac{\mu^+ P - 2F}{(1 - \mu^-) (R - P) + \mu^+ P} = q_2$$

Thus,  $q_1 < q < q_2$  with

$$F < \frac{P(R - P)(\mu^+ + \mu^- - 1)}{2R} = F_0$$

That is, the detection cost have to be lower than the limit ( $F_0$ ) since otherwise agents would never use this detection technology. We will assume that this condition is checked. Remark that  $P/R \in [q_1, q_2]$  and that  $0 < A_3 < q_1 < P/R < q_2 < A_4 < 1$ .

We postulate a linear relationship for the dynamic process  $\dot{l}_t = k [EG_I(q_t) - EG_U(q_t)]$  where  $k$  is a positive constant. Thus  $\dot{l}_t > 0$  if  $EG_I(q_t) - EG_U(q_t) > 0$  and we can write for the interval considered  $q \in (A_3, A_4)$  :

$$A_3 < q < q_1 \Rightarrow \dot{l}_t < 0 \quad (1)$$

$$q_1 < q < q_2 \Rightarrow \dot{l}_t > 0 \quad (2)$$

$$q_2 < q < A_4 \Rightarrow \dot{l}_t < 0 \quad (3)$$

Let us now analyze the dynamics of the moral agents ( $q$ ). The expected gain of a  $e_2^-$ - $U$ -type agent playing second is (with probability 1/2) :

$$\begin{aligned} EG_{e^-}(l) &= l [\mu^- R + (1 - \mu^-) P] + (1 - l)R \text{ for } q > P/R \\ &= l [\mu^- R + (1 - \mu^-) P] + (1 - l)P \text{ for } q < P/R \end{aligned}$$

The expected gain of a  $e_2^+$ - $U$ -type agent playing second is (with probability 1/2) :

$$\begin{aligned} EG_{e^+}(l) &= l [\mu^+ P + (1 - \mu^+) T] + (1 - l)T \text{ for } q > P/R \\ &= l [\mu^+ P + (1 - \mu^+) T] + (1 - l)P \text{ for } q < P/R \end{aligned}$$

For  $q > P/R$  we have :

$$EG_{e^-}(l) > EG_{e^+}(l) \iff l > \frac{(T - R)}{\mu^+ (T - P) - (1 - \mu^-) (R - P)} = l_0$$

The proportion  $l_0$  is positive as long as the denominator is, *i.e.* if the following condition holds :

$$\frac{T - P}{R - P} > \frac{1 - \mu^-}{\mu^+}$$

For  $q < P/R$  we obtain the reversed relation.

Let us further assume  $\dot{q}_t = h [EG_{e^-}(l_t) - EG_{e^+}(l_t)]$  with  $h$  a positive constant. Hence, we have  $\dot{q}_t > 0$  for  $EG_{e^-}(l_t) > EG_{e^+}(l_t)$ .

For  $\mu^- (R - P) > (1 - \mu^+) (T - P)$  we have :

$$\begin{aligned} 0 &< l < l_0 \text{ and } q < P/R \Rightarrow \dot{q}_t > 0 \\ 0 &< l < l_0 \text{ and } q > P/R \Rightarrow \dot{q}_t < 0 \\ l_0 &< l < 1 \Rightarrow \dot{q}_t > 0 \end{aligned}$$

Plotting this information and the information given by equations (1)-(3) we get figure (8).

For  $\mu^- (R - P) < (1 - \mu^+) (T - P)$  we have :

$$\begin{aligned} 0 &< l < l_0 \Rightarrow \dot{q}_t < 0 \\ l_0 &< l < 1 \text{ and } q > P/R \Rightarrow \dot{q}_t > 0 \\ l_0 &< l < 1 \text{ and } q < P/R \Rightarrow \dot{q}_t < 0 \end{aligned}$$

With this information and that of equations (1)-(3) we get figure (9).

In both cases, the unique rest point is given by the solution for  $(\dot{l}_t, \dot{q}_t) = (0, 0)$ . This rest point is  $(l^*, q^*) = (l_0, q_2)$ . This equilibrium implies a positive proportion of  $I$ -type agents and of  $e^-$ -type agents since  $l_0 > 0$  and  $q_2 > 0$ . However, this equilibrium is a vortex (see figures (8) and (9)) and has only the degenerate attraction set  $\{(p^*, l^*)\}$ , *i.e.* it will only be reached if the initial conditions are precisely the equilibrium. All other starting points, above  $P/R$  and below  $A_4$ , will lead to an indefinite cycling around  $(l^*, q^*)$ .